

PLAINTIFFS' EXHIBIT 74

**REDACTED VERSION OF DOCUMENT SOUGHT
TO BE FILED UNDER SEAL**

GIBSON DUNN

Gibson, Dunn & Crutcher LLP
555 Mission Street
San Francisco, CA 94105-0921
Tel 415.393.8200
www.gibsondunn.com

Rosemarie T. Ring
Direct: +1 415.393.8247
Fax: +1 415.801.7358
RRing@gibsondunn.com

HIGHLY CONFIDENTIAL—ATTORNEYS EYES ONLY

April 18, 2022

VIA JAMS ACCESS

Special Master Daniel B. Garrie
DGarrie@jamsadr.com

Re: *In re Facebook Consumer Privacy User Profile Litigation*, JAMS Ref No.
1200058674

Dear Special Master Garrie,

Pursuant to the Special Master's Order Following March 9, 2022 Hearing Regarding Plaintiffs' Motion to Compel Production of Plaintiff Data, dated March 22, 2022, Facebook submits this proposal to assist in determining "what, if any" additional named plaintiff data should be produced, as directed in Judge Corley's order dated January 12, 2022 (Dkt. No. 807), and respond to "Scenario" questions.

Facebook's Proposal For Producing Additional Named Plaintiff Data

As Judge Corley explained in her most recent order on named plaintiff data, issued on January 12, 2022, the Special Master is working with the parties to determine "what, if any" additional named plaintiff data should be produced consistent with Rule 26. Facebook continues to believe that the named plaintiff data that has already been produced, which includes data in all three categories of "discoverable user data" identified in Judge Corley's Discovery Order No. 9, satisfies its obligations under Rule 26. The additional data Plaintiffs seek is neither relevant to their claims nor proportional to the needs of this case because it was not shared or made accessible to third parties.

That said, in an effort to resolve this issue, and consistent with Judge Corley's Discovery Order Nos. 11 and 12 allowing Plaintiffs to test Facebook's position on whether there is additional named plaintiff data that was shared or made accessible to third parties, Facebook proposes the protocol set forth below for producing additional named plaintiff data.

As explained in Facebook's March 7, 2022 submission to the Special Master, at a high level, user data is stored in production systems and the data warehouse. [REDACTED]. Accordingly, we propose to produce additional data from these systems as explained below.

GIBSON DUNN

Special Master Daniel B. Garrie
 April 18, 2022
 Page 2

TAO. As also explained in the March 7, 2022 submission, TAO is a distributed data store for the social graph. There is an object and associations to that object for all users. *See* Ex. A, Internal Facebook Wiki Regarding TAO Core Concepts. For each named plaintiff, Facebook will produce the user objects and associations to those objects.

Hive. As explained in prior submissions to the Special Master, Hive is Facebook's data warehouse. *See* Ex. B, Declaration of Menggee Ji In Support of Facebook Inc.'s Motion For Reconsideration of Special Master's Order Regarding Named Plaintiff Data, ¶ 11.

Id. ¶ 16.

Id. ¶ 17.

Id. ¶ 18.

Id. These searches are subject to many technical and policy limitations which make searching for individual user data unduly burdensome, including that (1)

(*id.* ¶ 20); (2)

(*id.* ¶ 21); (3)

(*id.* ¶ 22); and (4) data in many Hive tables are in cold storage and would have to be restored in order to be searched and analyzed, *see* Ex. C, Declaration of Menggee Ji In Support of Facebook Inc.'s Motion For a Protective Order Against Production of API Call Logs, ¶¶ 17-18; Special Master's Nov. 8, 2021 Order Re: Facebook's Motion For Protective Order Against Production of API Call Logs, ¶ 15 ("Special Master Garrie finds that the data in the Mobile Table and Web Table is not reasonably accessible because it is not readily usable in its 'cold storage' state and must be restored to 'warm storage' in order to be searched and analyzed (i.e. usable).").

For these reasons, among others, as reported in its April 11, 2022 submission, Facebook estimates that it would take approximately [REDACTED] or around [REDACTED] of computational cost to search across the entire Hive data warehouse and extract all data about a single user.

Once the searches have been conducted, and data is returned, the data must be manually reviewed to, among other things, confirm that the data is associated with a particular user and does not include personal data of other users (e.g., user blocking another user) or trade secrets, and does not create system integrity or security concerns.

GIBSON DUNN

Special Master Daniel B. Garrie
April 18, 2022
Page 3

In light of the above technical limitations and burdens, Facebook proposes to search for and produce the following named plaintiff data from Hive to the extent it exists:

(1) specific types of data requested or referenced by Plaintiffs in challenging Facebook's production of named plaintiff data; (2) a sample of data from Hive tables with user identifiers included in Exhibit B to the April 11, 2022 submission; and (3) the only data identified in Exhibit C to the April 11, 2022 submission that has not already been produced to Plaintiffs in the DYI file or otherwise: set permissions (audience controls on a post).

1. Specific Types of Data

- *Any remaining off-platform activity:* DYI includes off-platform activity. Facebook will search for and produce any underlying raw log-level data for off-platform activity provided to Facebook by a third party associated with the named plaintiffs.
- *Any remaining ad interests:* DYI includes ad interests. Facebook will search for and produce any underlying raw log-level data associated with the named plaintiffs.
- *Any remaining ad click data:* DYI includes ad click data. Facebook will search for and produce any underlying raw log-level ad clicks data associated with the named plaintiffs.
- *Ad impressions data:* Facebook will search for and produce ad impressions data associated with the named plaintiffs.
- *Any remaining custom audience data:* DYI identifies third parties who have created custom audiences associated with a user. Facebook will search for and produce any more granular information about custom audiences associated with the named plaintiffs, including custom audience type, whether it was used to deliver ads, and when.

2. Sampling from Hive tables identified in Exhibit A of April 11, 2022 submission

Given that data in Hive tables are not shared or made accessible to third parties, the significant burden of searching for and producing individual user data from Hive as explained in prior submissions, and that Facebook has already produced nearly a million pages of the named plaintiff data from the DYI system, Hive data is not relevant or proportional to the needs of this case. Facebook nevertheless understands that, as Judge Corley found in Discovery Order Nos. 11 and 12, Plaintiffs are allowed to test Facebook's positions on "sharing and accessibility" by obtaining discovery into what data exists and how it is used. To that end, and consistent with the requirements of relevance and proportionality under Rule 26, Facebook proposes a sampling process.

GIBSON DUNN

Special Master Daniel B. Garrie
 April 18, 2022
 Page 4

In addition to the above specific categories of data, Facebook will produce a sample of data in warm storage from Hive tables identified in Exhibit A of the April 11, 2022 submission. Specifically, Facebook will produce data from 200 tables in Exhibit A to be jointly selected by Facebook and Plaintiffs. Facebook will randomly select 100 tables and Plaintiffs will select 100 tables.

3. User data identified in Exhibit B of the April 11, 2022 submission and not in DYI

As shown in Exhibit C of Facebook's April 11, 2022 submission to the Special Master, only two types of user data in the contracts submitted to the Special Master for in camera review are not included in the DYI file: (1) privacy settings, and (2) set permissions (audience controls on a post). (1) has already been produced. Facebook will agree to produce (2) from TAO.

With respect to other Facebook systems, not including TAO and Hive, there are seven systems¹ that did not fall into one of the categories identified by the Special Master in Hearing Order Regarding Plaintiffs' Motion To Compel Production of Plaintiff Data, dated February 21, 2022: [REDACTED], [REDACTED] (MySQL), [REDACTED] (MySQL), [REDACTED], [REDACTED], and [REDACTED]. Producing data from any of these systems would be extremely burdensome and is not relevant or proportional to the needs of this case.

- [REDACTED] and [REDACTED]. These systems are BLOB and binary storage and are not searchable by user identifiers.
- [REDACTED]. This system is a real-time aggregation counter and is not searchable by user identifiers.
- [REDACTED] (MySQL), [REDACTED] (MySQL), [REDACTED] and [REDACTED]. These systems contain unstructured data and cannot be searched by user identifiers, except for tables using Ent or Node schemas for which there is a structured deletion plan. Searching other tables for named plaintiff data would require a full text search of all tables that would take many months to complete and generate a significant number of false positives requiring manual review of all data returned.

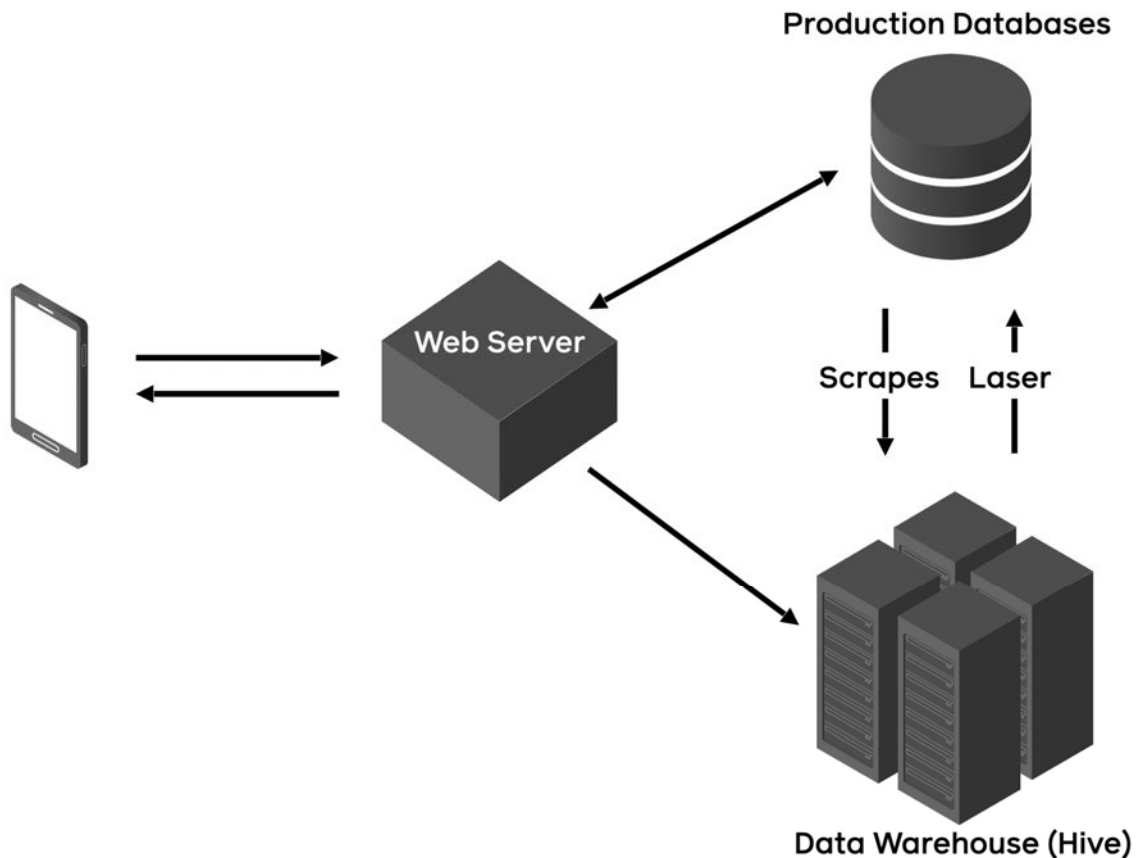
With respect to Dataswarm, the Special Master suggested that Facebook "query Tasks to identify Task Definitions which involve named plaintiff data" using identifiers,

¹ [REDACTED] should have been put in Category (6) because the data is in DYI. [REDACTED] should also have been put in Category (6) because it is duplicative of TAO. MySQL describes the general use of MySQL databases and should not be a separate entry. For example, [REDACTED] (MySQL) and [REDACTED] (MySQL) are specific uses of MySQL.

GIBSON DUNN

Special Master Daniel B. Garrie
 April 18, 2022
 Page 5

review those Tasks Definitions, and then search for named plaintiff data in the sources/destinations identified. We have investigated the Special Master's suggestion and did not include it in our proposed protocol because it is not feasible. Dataswarm is a collection of data processing operations coded in Python. Task Definitions are the Python code and Tasks are the running of that code. Searching Python code for user identifiers would be done through a syntax-only search which cannot understand the semantic meaning of the data being processed, i.e., whether a column name being processed involves a user identifier. As a result, additional code would have to be developed to understand whether/how identifiers are involved in a particular Task Definition.



GIBSON DUNN

Special Master Daniel B. Garrie
April 18, 2022
Page 6

Scenarios

No later than April 4, 2022, Facebook is to submit documentation sufficient to describe the data collected both on and off platform or provided by Third Parties in the following scenarios and provide written responses to the questions below: “Exhibit A to Plaintiff’s Questions re: Data Collection and Use indicates that Facebook used predictive algorithms to generate five political segments for Facebook users (Very Liberal, Liberal, Moderate, Conservative, and Very Conservative) based on demographic, psychographic, and behavioral signals from Facebook user data.”

1. What are the inputs into these algorithms (i.e. what are the demographic, psychographic, and behavioral signals used to generate the political segments)? Are these inputs provided by users or derived by Facebook?

Answer: As an initial matter, Very Liberal, Liberal, Moderate, Conservative, and Very Conservative are ad interests (“political ad interests”), which are distinct from political segmentation “based on demographic, psychographic, and behavioral signals from Facebook user data” (“political segments”). Political segments are political ad interests that are divided into political segments.

Political ad targeting segments were deprecated for new campaigns in January 2022 and for existing campaigns in March 2022. Based on a reasonable investigation to date, these were ad targeting segments derived from the following on-platform activity (user-provided data and observed data):

[REDACTED]

Political segments were deprecated for new campaigns in October 2017 and for existing campaigns in January 2018. Based on a reasonable investigation to date, these segments were created by

[REDACTED]

GIBSON DUNN

Special Master Daniel B. Garrie
April 18, 2022
Page 7

2. How are the psychographic signals computed (e.g. how is the psychographic signal “High Dollar Religious Donor” determined)?

Answer: Based on a reasonable investigation to date, we believe “High Dollar Religious Donor” was an ad targeting option based on partner categories. We do not know how this partner category was created. Partner categories were deprecated in October 2018.

3. Is information regarding identifiable ethnic affinities provided by users or derived by Facebook? How is ethnic affinity derived?

Answer: Ethnic affinities, also referred to as multicultural affinities, were ad interests created by Facebook based on a user’s on-platform activity indicating an interest in content relating to certain communities. Ethnic affinity ad interests were deprecated in August 2020.

4. Where is political segmentation data for Facebook users stored?

Answer: Political segments were deprecated in October 2017 for new campaigns. Existing campaigns using political segments ended on January 1, 2018. Based on a reasonable investigation, we believe that political segments were deleted in January 2018.

5. Is political segmentation determined for a Facebook user as part of a data process that runs on a regularly scheduled basis or evaluated in real time when an ad is served?

Answer: Based on a reasonable investigation to date, we believe political segments were created manually by [REDACTED]

6. Is political segmentation associated with a Facebook user if possible (i.e. via UID, RID, SID, ASID, or another identifier that can be mapped to a user)? If so, explain how the political segmentation is associated with a Facebook user.

Answer: Based on a reasonable investigation to date, we believe political segments were associated with users by UIDs.

7. Is an individual’s assigned political segment part of the DYI file?

Answer: Based on a reasonable investigation to date, we do not believe political segments were included in the DYI file.

GIBSON DUNN

Special Master Daniel B. Garrie
April 18, 2022
Page 8

Sincerely,

A handwritten signature in cursive script that reads "Rosemarie Ring". The ink is dark and the signature is fluid.

Rosemarie T. Ring

EXHIBIT A

EXHIBIT B

GIBSON, DUNN & CRUTCHER LLP
Orin Snyder (*pro hac vice*)
osnyder@gibsondunn.com
200 Park Avenue
New York, NY 10166-0193
Telephone: 212.351.4000
Facsimile: 212.351.4035

Kristin A. Linsley (SBN 154148)
klinsley@gibsondunn.com
Martie Kutscher (SBN 302650)
mkutscherclark@gibsondunn.com
555 Mission Street, Suite 3000
San Francisco, CA 94105-0921
Telephone: 415.393.8200
Facsimile: 415.393.8306

GIBSON, DUNN & CRUTCHER LLP
Deborah Stein (SBN 224570)
dstein@gibsondunn.com
333 South Grand Avenue
Los Angeles, CA 90071-3197
Telephone: 213.229.7000
Facsimile: 213.229.7520

Joshua S. Lipshutz (SBN 242557)
jlipshutz@gibsondunn.com
1050 Connecticut Avenue, N.W.
Washington, DC 20036-5306
Telephone: 202.955.8500
Facsimile: 202.467.0539

Attorneys for Defendant Facebook, Inc.

**UNITED STATES DISTRICT COURT
NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO DIVISION**

IN RE: FACEBOOK, INC. CONSUMER
PRIVACY USER PROFILE LITIGATION,

This document relates to:

ALL ACTIONS

CASE NO. 3:18-MD-02843-VC

**DECLARATION OF MENGGE JI IN SUPPORT OF FACEBOOK, INC.'S
MOTION FOR RECONSIDERATION OF
THE SPECIAL MASTER'S ORDER RE
PLAINTIFFS' MOTION TO COMPEL
PRODUCTION OF PLAINTIFF DATA**

Discovery Special Master Daniel Garrie, Esq.

**DECLARATION OF MENGGE JI IN SUPPORT OF FACEBOOK, INC.'S
MOTION FOR RECONSIDERATION OF THE SPECIAL MASTER'S ORDER RE
PLAINTIFFS' MOTION TO COMPEL PRODUCTION OF PLAINTIFF DATA**

I, Mengge Ji, declare as follows:

1. I am a Data Scientist at Meta, Inc. f/k/a Facebook, Inc. (“Facebook”). My job responsibilities include, among other things, understanding and working with Facebook’s data systems, writing queries and conducting analyses of these data, researching Facebook’s data and related technologies, and locating, analyzing, and exporting data for production in litigation and other legal matters. I submit this declaration in support of Facebook’s motion for reconsideration of the Special Master’s Order re Plaintiffs’ Motion to Compel Production of Plaintiff Data. Unless otherwise stated, I have personal knowledge of the facts set forth herein, and, if called as a witness, I could and would competently testify thereto.

2. In my role as Data Scientist at Facebook, I am familiar with the general categories of user information Facebook maintains, where that data might be stored, and how they can potentially be accessed.

Data Contained in Facebook’s Social Graph Related to Users’ Activities on Facebook

3. Facebook users directly provide Facebook with data or information when using the Facebook platform. For instance, a user can upload a photo directly to her Facebook profile, which Facebook retains in order to display when that profile is accessed (if the user’s privacy settings allow it). If the user “tags” a friend in that photo, Facebook also retains a record of which other Facebook user has been identified as appearing in that photo. In essence, the entire Facebook product that users interact with is a web of data points and relationships between data points that Facebook’s systems present in a user-friendly format.

4. Facebook uses the term the “Social Graph” to describe this complex web of people, places, things, actions, and connections on the Facebook platform. As Facebook users navigate through Facebook and interact with it—including, for example, by commenting on

posts made by other users, watching videos, posting photos, and sending messages—the users create new relationships and connections between themselves and the content they are able to see. The activities a user takes on the Facebook platform, including posting pictures to their profiles, tagging friends in photos, and commenting on friends’ timelines, are reflected in the Social Graph.

5. The Facebook product that users see is powered by a series of databases that work in tandem to provide Facebook users a seamless experience.¹ The key databases Facebook uses to support the Facebook product, which store the user content and information presently accessible via the Social Graph, are:

- [REDACTED]
- [REDACTED]
- [REDACTED]
- [REDACTED]

¹ Facebook’s Social Graph is not powered by a single database. Rather, Facebook is powered by an extraordinarily complex information architecture that stores information in various databases. The information in these databases is generally not human-readable and instead is intended to be processed for human consumption through Facebook’s production environment.

- [REDACTED]
[REDACTED]
[REDACTED].²

6. Between them, these databases contain trillions of data points provided to Facebook by its more than 2 billion users, many of which need to be readily accessible at all times so that users can view them as they use the Facebook product. As a result, the data infrastructure underlying the Social Graph is not only extraordinarily complex, but has also been specifically designed to serve the needs of the Facebook product.

7. One of the functions of the Facebook product is to be able to specifically identify user profiles or other sets of information associated with specific users. [REDACTED]

[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]

8. [REDACTED]

[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]

9. I understand that Facebook’s “Download Your Information” or “DYI” tool retrieves data from and allows users to download a copy of data Facebook associates with their

² The Everstore database is in the process of being deprecated. The data contained in Everstore is being transitioned to Manifold.

Facebook account, including data associated with their account in the Social Graph. In order to compile this information, the DYI tool runs a set of searches that have been engineered by the DYI product team in order to pull data associated with that user from the Social Graph. The resulting data is presented in a user-readable format, rather than as code or strings of data.

Analytics Data Stored in Facebook's Data Warehouse ("Hive")

10. Facebook also stores data sets that it uses for internal analytics, product development, and other business functions. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

11. Facebook's internal analytics data is primarily stored in a data warehouse called Hive, which exists separately from the Social Graph. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

12. The Hive data warehouse is contained across a decentralized set of data centers, spread across the world. [REDACTED]

[REDACTED]

13. In addition to the Social Graph and Hive, I am aware of several other databases that contain data related to Facebook users, including data that is anonymized. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

14. I have also consulted with other members of the Data Science team regarding this question. [REDACTED]

[REDACTED]

15. [REDACTED]

[REDACTED]

[REDACTED]

Searching Within Hive

16. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

17. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

18. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

19. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

20. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

21. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

22. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

23. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

Exporting Data from Hive

24. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

25. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED].

26. [REDACTED]

[REDACTED]

27. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

[REDACTED]

28. [REDACTED]

[REDACTED]

[REDACTED]

29. Preparing any productions also requires substantial additional resources from senior members of the eDiscovery data scientist team to prepare for and manage the coding and export project. Moreover, after the initial exported files have been created, it would take additional time for a review team to validate the queries and data and perform additional quality control, or longer if there are issues that require remediation.

30. [REDACTED]

[REDACTED]

[REDACTED]

[REDACTED] The Facebook eDiscovery team, which is responsible for assisting Facebook in-house and outside counsel in active litigations and other legal matters, in addition to building and maintaining internal infrastructure crucial to the management and preservation of data on legal hold, does not have the time and resources required to search for, access, analyze, and export the data in millions of Hive tables in the manner described above.

31. Facebook cannot materially shorten this timeline by hiring new employees

[REDACTED]

[REDACTED]

[REDACTED] For the same reason, Facebook cannot simply engage third party consultants or temporary employees to handle this data export. Nor would adding more servers—which would require diverting them from their use in the ordinary course of business—necessarily reduce the estimated timeline in a linear fashion, [REDACTED]

[REDACTED]

[REDACTED] All of these options—hiring new employees, hiring contractors, and adding servers—would also be extremely costly.

I declare under penalty of perjury that the foregoing is true and correct, and that I executed this Declaration on December 10, 2021, in Sausalito, California.



Mengge Ji

EXHIBIT C

GIBSON, DUNN & CRUTCHER LLP
Orin Snyder (*pro hac vice*)
osnyder@gibsondunn.com
200 Park Avenue
New York, NY 10166-0193
Telephone: 212.351.4000
Facsimile: 212.351.4035

Kristin A. Linsley (SBN 154148)
klinsley@gibsondunn.com
Martie Kutscher (SBN 302650)
mkutscherclark@gibsondunn.com
555 Mission Street, Suite 3000
San Francisco, CA 94105-0921
Telephone: 415.393.8200
Facsimile: 415.393.8306

GIBSON, DUNN & CRUTCHER LLP
Deborah Stein (SBN 224570)
dstein@gibsondunn.com
333 South Grand Avenue
Los Angeles, CA 90071-3197
Telephone: 213.229.7000
Facsimile: 213.229.7520

Joshua S. Lipshutz (SBN 242557)
jlipshutz@gibsondunn.com
1050 Connecticut Avenue, N.W.
Washington, DC 20036-5306
Telephone: 202.955.8500
Facsimile: 202.467.0539

Attorneys for Defendant Facebook, Inc.

**UNITED STATES DISTRICT COURT
NORTHERN DISTRICT OF CALIFORNIA
SAN FRANCISCO DIVISION**

IN RE: FACEBOOK, INC. CONSUMER
PRIVACY USER PROFILE LITIGATION,

This document relates to:

ALL ACTIONS

CASE NO. 3:18-MD-02843-VC

**DECLARATION OF MENGGE JI IN
SUPPORT OF FACEBOOK, INC.'S
MOTION FOR A PROTECTIVE ORDER
AGAINST PRODUCTION OF API CALL
LOGS**

Discovery Special Master Daniel Garrie, Esq.

**DECLARATION OF MENGGE JI IN SUPPORT OF
FACEBOOK, INC.'S MOTION FOR A PROTECTIVE ORDER AGAINST
PRODUCTION OF API CALL LOGS**

HIGHLY CONFIDENTIAL

I, Mengge Ji, declare as follows:

1. I am a Data Scientist at Facebook, Inc. (“Facebook”). My job responsibilities include, among other things, understanding and working with Facebook’s data systems, writing queries and conducting analyses of these data, researching Facebook’s data and related technologies, and locating, analyzing, and exporting data for production in litigation and other legal matters. I submit this declaration in support of Facebook’s motion for a protective order as against the production of certain of Facebook’s API call logs. Unless otherwise stated, I have personal knowledge of the facts set forth herein, and, if called as a witness, I could and would competently testify thereto.

2. My understanding is that Plaintiffs in this matter have requested that Facebook produce to them all “logs of API Calls made by Third Parties for the Content and Information of Friends of Installing Users,” including specific information about, among other things, “the response status of those Calls (*e.g.*, success or reason for failure)” and “[t]he Content and Information delivered to Third Parties on each API in response to these Calls.”

3. I also understand that Plaintiffs have issued additional requests for information relating to the number of API calls made by certain “Third Parties” and “Business Partner[s],” including “the number of calls . . . made to each such API . . . each month” and “the volume of data transferred from each such API . . . each month.”

1. Background on Facebook’s APIs and Call Logs

4. Application Program Interfaces (“APIs”) are standard computing protocols used throughout the digital world. Each time a user visits an app, the user necessarily interacts through an API. And every company that allows third parties to communicate with their servers

uses APIs. As just one example, email programs like Outlook or Mail on iOS communicate with the user's email provider (Yahoo, Gmail, etc.) through APIs to obtain the user's messages.

5. Facebook has a series of APIs that allow for data to be transferred to and from Facebook's Platform. The Facebook Platform is powered by a series of databases that Facebook relies upon to host and run the Facebook products and which work in tandem to provide Facebook users a seamless experience as they access different categories and types of content posted on the Platform, by themselves or other users. APIs facilitate the transfer of data from the databases that support the Platform to websites and applications, including Facebook's own, proprietary mobile apps. For instance, if a Facebook user pulls up her profile photo via the Facebook mobile app, the Facebook mobile app will "call" a Facebook API to request that photo to be displayed locally on the user's phone.

6. Third party mobile apps and websites can also connect to Facebook through its APIs. For instance, if a user chooses to log in to NYTimes.com with his Facebook account instead of creating an account for that site, NYTimes.com will call the Facebook APIs to verify the user's identity. Facebook has a series of public APIs that are published on its developer website (<https://developers.facebook.com/>), and app developers can use them or request to use them to send or receive information from the Facebook Platform.

7. Facebook maintains internal logs that contain the number of API calls made to the Facebook Platform. Two of the data tables that contain [REDACTED] information Facebook maintains regarding API call activity are the [REDACTED] table and the [REDACTED] table.

2. The [REDACTED] and [REDACTED] Tables

8. The [REDACTED] table contains logs of API calls [REDACTED]
[REDACTED]

[REDACTED]

[REDACTED] The data in the [REDACTED] table dates back to [REDACTED]. As of October 6, 2021, Facebook had [REDACTED] of data in this table.

9. The categories of data logged in the [REDACTED] table include, among other things, [REDACTED]. A true and correct list of the data fields reflected in the [REDACTED] table is attached hereto as **Exhibit A**. A true and correct sample of [REDACTED] from the [REDACTED] table, with [REDACTED] redacted, is attached hereto as **Exhibit B**.

10. The [REDACTED] table contains logs of API calls [REDACTED]. The data in the [REDACTED] table dates back to [REDACTED]. As of October 6, 2021, Facebook had [REDACTED] days of data in this table.

11. The categories of data logged in the [REDACTED] table include, among other things, [REDACTED]. A true and correct list of the data fields reflected in the [REDACTED] table is attached hereto as **Exhibit C**. A true and correct sample of ten lines of data from the [REDACTED] table, with the [REDACTED] redacted, is attached hereto as **Exhibit D**.

12. Several of the fields in each of these tables do not contain [REDACTED]

[REDACTED]

[REDACTED]. This is an extra step of analysis that would need to be conducted if these tables were ever to be used or analyzed outside of Facebook's environment.

13. The [REDACTED] and [REDACTED] tables do not contain data demonstrating whether data was returned in response to an API call, what the data that was returned may have been (if it was returned at all), or the volume of data that may have been returned (if it was transferred at all).

14. I understand that there are many reasons why data would not be returned in response to an API call. These reasons include, but are not limited to, [REDACTED]

[REDACTED]
[REDACTED]
[REDACTED]

A. Storage of the [REDACTED] and [REDACTED] tables

15. The [REDACTED] and [REDACTED] tables are stored in a highly compressed state in Facebook's internal decentralized, structured data warehouse, known as "Hive." Hive is designed to ingest and process data for internal analysis, not to export data into files that can be transferred outside of Facebook. As discussed further below, accessing and analyzing large data tables in Hive is a multi-step process that can take a considerable amount of time, depending on the volume of data being accessed and the complexity of the analysis or querying needed.

16. The active [REDACTED] and [REDACTED] tables are currently maintained pursuant to [REDACTED] business retention periods [REDACTED]

[REDACTED]
[REDACTED]

[REDACTED]

[REDACTED]

17. [REDACTED]

[REDACTED]

[REDACTED]

18. Tables in cold storage cannot be reviewed or analyzed while they remain in cold storage. If Facebook wants to conduct any analysis of data in cold storage, a Facebook data scientist must first restore it from cold storage to an accessible location, referred to as “warm” storage. It takes approximately 1 day of processing to restore [REDACTED] of cold storage data. Restored cold storage data is only made available to legal team members.

19. The [REDACTED] and [REDACTED] tables are stored in “partitions” or segments. These tables are partitioned [REDACTED]. Because each [REDACTED] is partitioned off from the others, Facebook cannot quickly or efficiently run searches across multiple [REDACTED]. Instead, Facebook’s system is optimized for running searches within an individual partition. As a result, Facebook cannot search these entire tables for a specific value, for instance a [REDACTED], while in cold storage, and instead must restore individual partitions to warm storage and then search the partitions individually. Due to the same limitations on accessing and searching multiple partitions of data, Facebook also cannot filter these tables for [REDACTED] [REDACTED] without first restoring partitions from cold storage piecemeal.

20. Accordingly, to search or analyze the data in either the [REDACTED] table or [REDACTED] tables, a Facebook data scientist must first restore a partition of the data to warm

storage. Due to limitations on computing and server capacity, only a limited number of partitions can be restored and analyzed at any given time.

21. The [REDACTED] and [REDACTED] tables are very large data tables, even in comparison to other data sets at Facebook and even when compressed in Facebook's specialized data storage environment. The volume of the [REDACTED] table is currently [REDACTED] in its most highly-compressed form in cold storage. The volume of the [REDACTED] table is currently [REDACTED] in its most highly-compressed form. Maintaining the high volume of this data [REDACTED] is beginning to tax Facebook's data storage system, which is one of the largest and most sophisticated data storage systems in the world. For instance, I understand that it is costing Facebook approximately [REDACTED] to store the [REDACTED] data in its compressed form.

B. Sampling of the [REDACTED] table

22. In order to assess how long it might take to produce data from the [REDACTED] table to the Plaintiffs in this action, I recently restored and analyzed one partition of data from the [REDACTED] table dated [REDACTED].

23. The volume of the [REDACTED] partition is [REDACTED] in cold storage. For reference, 1 TB is four-times the size of a standard laptop computer hard drive of 256 GB, so the *compressed* volume of the data for this one day equals [REDACTED]. The size of the data in each partition in the [REDACTED] table varies day [REDACTED] [REDACTED]. For example, the most recent partitions of the table in cold storage are around [REDACTED] each.

24. In order to analyze the [REDACTED] data, I first restored this partition to warm storage. The length of time it takes to restore data from cold storage depends on how many

requests are pending at the same time and the computing power available, which varies depending on the time of day. In this instance, it took approximately 2 hours of processing time to restore this single partition into warm storage. Based on the volume of this partition and others in this table, I estimate that Facebook presently has the computing and server capacity to host up to five partitions of this table in warm storage at a time. As a result, Facebook can analyze no more than five partitions of this table at a given time in the ordinary course. In other words, an effort to simultaneously analyze more than five partitions of this data would displace and disrupt all other efforts by the legal team to perform other data analyses for this and other matters, including for the [REDACTED] table.

25. The [REDACTED] partition of the [REDACTED] table contains [REDACTED]. As a point of reference, a Microsoft Excel file can only hold up to 1 million rows of data. If translated into Excel, then, this one partition of the [REDACTED] table would comprise more than [REDACTED] Excel files.

26. As I mentioned in paragraph 8 above, this table only contains logging for API calls [REDACTED]
[REDACTED]
The partition only contained data for [REDACTED]. It took 1 hour and 14 minutes for this search to complete.

27. [REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]

28. Based on the results of this exercise, I estimate it would take 2,378 hours ([REDACTED] in cold storage) to restore each of the partitions in the [REDACTED] table to warm storage and additional 1,279 hours ([REDACTED]) to run a single search on each. Assuming I worked exclusively on this task for 8-hours per day, 260 days each year, it would take nearly two years for me to run this same search on [REDACTED] of this data (through the [REDACTED] partition). This is also a conservative estimate, [REDACTED]. This estimate also does not account for the other demands on Facebook's personnel and resources (including the need to use these same systems to respond to requests for other legal matters), occasional outages, or many other variables.

C. Sampling of the [REDACTED] table

29. I also restored and recently analyzed one partition of data from the [REDACTED] table dated [REDACTED].

30. The volume of the [REDACTED] partition is [REDACTED] in cold storage. The size of the data in each partition in the [REDACTED] table varies [REDACTED]. For example, the most recent partitions of the table in cold storage are around [REDACTED] each.

31. In order to analyze the [REDACTED] data, I first restored this partition to warm storage. The length of time it takes to restore data from cold storage depends on how many requests are pending at the same time and the computing power available, which varies depending on the time of day. In this instance, it took approximately [REDACTED] hours of processing time to restore this partition into warm storage for analysis. Based on the volume of this partition and others in this table, I estimate that Facebook presently has the computing and server

capacity to host 20 partitions of this table in warm storage at a time. As a result, Facebook can only analyze 20 partitions of this table at a given time in the ordinary course.

32. The [REDACTED] partition of the [REDACTED] table contains

[REDACTED] rows of information. If exported to Excel, then, this partition of the

[REDACTED] table would alone comprise [REDACTED] Excel files.

33. I ran a search for [REDACTED]

[REDACTED]. The partition only contained data for [REDACTED].

It took 43 minutes for this search to run.

34. Of the [REDACTED] of data in this table, [REDACTED] of data were

returned for these [REDACTED]. [REDACTED]

[REDACTED]

35. Based on the results of this exercise, I estimate it would take 2,970 hours (2.5

hours per partition in cold storage) to restore each of the partitions in the [REDACTED] table to

warm storage and an additional 935 hours (45 minutes per partition) to run a single search on

each. Assuming I worked exclusively on this task for 8-hours per day, 260 days each year, it

would take nearly two years for me to run this same search on each partition of this data (through

the [REDACTED] partition). This is also a conservative estimate, [REDACTED]

[REDACTED]. This estimate also does not account

for the other demands on Facebook's personnel and resources, occasional outages, or many other

variables. This estimate is also *in addition to* the two years it would take to search the

[REDACTED] table.

36. As a result, it is neither feasible nor practicable to restore all of the partitions of the [REDACTED] and/or [REDACTED] tables at the same time and [REDACTED]

3. The Burden of Production of API Call Log Data

37. As noted, the [REDACTED] and [REDACTED] tables are stored in Facebook's internal structured data warehouse, known as "Hive." This data warehouse is not optimized to export data in a manner that can be transferred outside of Facebook. For Facebook to export and produce data from Hive tables, the data must be decompressed, broken into small pieces, and re-written to .csv files (which can be opened with Excel). This would be a manual, iterative process conducted by Facebook data scientists on the legal team. As a result, the timeline for completing an export needs to account not only for the limitations on Facebook's computing capacity, the size limitations on our mechanisms for export, and the availability of Facebook's technical resources, but also the possibility of human error, which necessitates real-time quality checks of the data downloads.

38. Before any analysis or export can be conducted, the data must be transitioned to warm storage in order to be reviewed, analyzed, or exported. As discussed above in Paragraphs 28 and 35, I estimate it would take 5,348 hours to restore all of the partitions of [REDACTED] and [REDACTED] tables in to warm storage, which is more than 2.5 years of work, assuming I work exclusively on this task for 8 hours per day, 260 days each year, and do not run any searches to narrow the data set. Running searches on the data set would take at least an additional year.

39. After data has been restored into warm storage, it is necessary to decompress the data, which requires that Facebook take several additional steps, which are themselves time-consuming.

40. First, a Facebook data scientist must write entirely new code to query the data and demarcate each partition into chunks of data that can be effectively transported to and/or received by the recipient. For instance, if the recipient intends to use Excel to analyze the data, each partition must be broken into chunks of 1 million rows each, which is the maximum number of rows that can be opened in an Excel file. When Facebook must produce Hive data in litigation, it typically produces it in Excel form. If the intended recipient—in this instance, Plaintiffs—had access to or intended to use specialized software to access the data (*e.g.*, SQL, R, Stata, SAS) and had sufficient computing power to access larger files, the limiting factor would be a combination of size of the drive used to transport the data and the maximum number of rows that could be reliably written into a .csv and reviewed for quality control. I understand from my colleagues at Facebook that the largest .csv files that we have handled and transferred reliably is 50 GB, which is roughly [REDACTED] [REDACTED] [REDACTED]

41. Second, the code is run to export that data to a server, which would involve querying each chunk of [REDACTED] from the data warehouse and then writing the results of that query as .csv files in a decompressed flat-file format. The speed at which Facebook can export the data from the warehouse to the server varies depending on usage of its systems in the ordinary course of business, but based on testing and prior exports, Facebook's current best estimate is that it takes about 40 minutes to write a 50 GB .csv file. At this rate, it would take approximately 15 hours to extract the [REDACTED] of data I identified in my

search of the [REDACTED] partition of the [REDACTED] table ([REDACTED]
[REDACTED]), discussed in Paragraph 27 above, to .csv format. The resulting .csv files would be approximately [REDACTED] total, or approximately [REDACTED]
[REDACTED]. It would then take our eDiscovery and Data Science teams about 3 minutes per file to inspect the data to confirm that it was not corrupted or truncated during export. If we were to export the [REDACTED] of data I identified in my search of the [REDACTED] partition of the [REDACTED] table—[REDACTED]
[REDACTED], it would take approximately 1.5 hours to inspect these files. It would therefore take approximately 16.5 hours to extract and inspect .csv files containing [REDACTED] [REDACTED] [REDACTED]

42. It would separately take 2 hours to extract the [REDACTED] data I identified in my search of [REDACTED] [REDACTED] discussed in Paragraph 34 above, to .csv format. The resulting .csv files would be approximately [REDACTED] total. It would then take our e-Discovery and Data Science teams about 3 minutes per file to inspect the data to confirm that it was not corrupted or truncated during export. If we were to export the [REDACTED] data I identified in my search of the [REDACTED] partition of the [REDACTED] [REDACTED]
[REDACTED], it would take approximately 12 minutes to inspect these files. It would therefore take an additional 2 to 3 hours to extract and inspect .csv files containing [REDACTED]
[REDACTED].

43. Third, once the server is full, a data scientist must empty the server by transferring the .csv files into a local encrypted hard drive or a secure transfer site. Each 50 GB .csv file

would require 5 to 8 hours to be transferred from the Facebook servers to a physical hard drive, depending on network speeds and other variables.

44. Decompression of the data, which is required for export of this data, results in a significant expansion in the size of the data. Within Facebook's systems, one partition of the [REDACTED] table comprises approximately [REDACTED] of data and one partition of the [REDACTED] table comprises approximately [REDACTED] of data. It is difficult to estimate the compression ratio before the data is actually exported. Based on testing and prior exports, Facebook's current best estimate as to how big a single partition of each of these tables would be when fully decompressed to .csv format for export is around [REDACTED] per partition for the [REDACTED] table and [REDACTED] per partition for the [REDACTED] table.

45. [REDACTED]
[REDACTED]. As a result, the second and third steps of this process would have to be run a minimum of 19,250 separate times for each partition of the [REDACTED] table. Similarly, the second and third steps of this process would have to be run a minimum of 4,750 separate times for each partition of the [REDACTED] table. It is not possible to expand the number of available servers without investing significantly in additional resources or interrupting other business and legal functions.

46. The decompressed volume of a single partition of data from each of these tables—up to [REDACTED] table and [REDACTED] for the [REDACTED] table—is significantly larger than any hard drives or secure transfer site that Facebook has access to. Facebook's current best estimate is that exporting a single partition of the [REDACTED] table would require it to load .csv files onto [REDACTED] encrypted hard drives that would cost approximately \$1,200 each, presuming Facebook can acquire a

sufficient number of them. Facebook's current best estimate is that exporting a single partition of the [REDACTED] table would require it to load .csv files onto [REDACTED] encrypted hard drives that would cost approximately \$1,200 each, again presuming Facebook can acquire a sufficient number of them.

47. The process must be repeated over and over until the export is complete. Facebook estimates that it would take 8,000 days and \$70,000 in hardware costs to export the smallest, single partition of the [REDACTED] table and 2,000 days and \$17,000 in hardware costs to export a single partition of the [REDACTED] table.¹ Replicating this process for more than [REDACTED] of each table would multiply those timing estimates proportionally.

48. This manual, iterative process introduces the possibility of human error, which necessitates real-time quality checks of the data downloads that take additional time and create the potential for delay.

49. Facebook also estimates that it will take several days for senior members of the eDiscovery data scientist team conducting the export to prepare for and manage the coding and export project. After the initial exported files have been created, Facebook estimates it would take an additional month for a review team to validate the queries and data and perform additional quality control, or longer if there are issues that require remediation.

A. Impact on Facebook's E-Discovery Team's Resources

50. In total, Facebook currently estimates that, even without any other responsibilities or deadlines, it would be substantial, multi-year project to attempt to export the [REDACTED] table and the [REDACTED] table through the process described above. The Facebook E-Discovery team, which is responsible for assisting Facebook in-house and outside counsel in

¹ The shipping costs for the hard drives would add additional costs onto this project.

active litigations and other legal matters, in addition to building and maintaining internal infrastructure crucial to the management and preservation of data on legal hold, does not have the time and resources required to access, analyze, and export the data in the [REDACTED] table and the [REDACTED] table in the manner described above.

51. Facebook cannot materially shorten this timeline by hiring new employees because it would take time and resources to interview, select, and onboard new employees, and any new hires would have to be sufficiently trained regarding Facebook's systems, policies, and procedures, which is a mandatory process that takes three weeks. The API call log data export in particular requires specialized knowledge about Facebook's systems because the system is not built for the purpose of exporting the extremely large amounts of data at issue—knowledge that non-Facebook personnel do not have, given the nature of the proprietary systems and databases employed by Facebook. For the same reason, Facebook cannot simply engage third party consultants or temporary employees to handle this data export. Nor would adding more servers—which would require diverting them from their use in the ordinary course of business—necessarily reduce the estimated timeline in a linear fashion, as the export still requires downloading the data from the server, and the manual process of writing the code and monitoring the export. All of these options—hiring new employees, hiring contractors, and adding servers—would also be extremely costly.

4. Alternative Production Proposal

52. I understand that Facebook has proposed to produce summary data relating to API calls to Plaintiffs in this matter instead of data from the [REDACTED] and [REDACTED] tables. [REDACTED]

[REDACTED]

[REDACTED].

53. Specifically, I understand that Facebook has proposed to produce data from the [REDACTED] table. The [REDACTED] table contains summaries of data [REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]
[REDACTED]

54. Because the [REDACTED] table summarizes [REDACTED], it is substantially less burdensome to store and analyze.

55. The [REDACTED] table contains summary data regarding API calls made to the Facebook Platform, [REDACTED]
[REDACTED]. The fields contained in this table include, among other things, [REDACTED]
[REDACTED]. A true and correct list of the data fields reflected in the [REDACTED] table is attached hereto as **Exhibit E**. A true and correct sample of [REDACTED] from the [REDACTED] table is attached hereto as **Exhibit F**.

56. The data in the [REDACTED] table dates back to [REDACTED].

57. The [REDACTED] table is maintained in warm storage in the Hive data warehouse. It is currently [REDACTED] in this compressed form. A single partition [REDACTED]

of this table is around [REDACTED] in compressed form and [REDACTED] in raw (*i.e.*, not compressed) form.

58. It is difficult to estimate the compression ratio before the data is actually exported. Based on testing and prior exports, Facebook's current best estimate as to how big an export of data from this table dating from [REDACTED] would be when fully decompressed to .csv format for export is [REDACTED].

59. Facebook is currently preparing an export of data from this table dating from [REDACTED] for production to Plaintiffs in this matter and estimates it will be able to produce this data by November 30, 2021.

60. A brief summary of the three tables discussed herein is provided below, along with relevant data points regarding the burden of production.

	[REDACTED]	[REDACTED]	[REDACTED]
Earliest Date	[REDACTED]	[REDACTED]	[REDACTED]
Compressed Size	[REDACTED]	[REDACTED]	[REDACTED]
Est. Decompressed Partition Size	[REDACTED]	[REDACTED]	[REDACTED]
Est. Decompressed Size of Data Through [REDACTED]	[REDACTED]	[REDACTED]	[REDACTED]

I declare under penalty of perjury that the foregoing is true and correct, and that I executed this Declaration on October 18, 2021, in Sausalito, California.

A handwritten signature in black ink, appearing to read "Mengge Ji", is written over a horizontal line.

Mengge Ji

Exhibit A

to Ji Declaration

Highly Confidential

Exhibit B

to Ji Declaration

Highly Confidential

Exhibit C

to Ji Declaration

Highly Confidential

Exhibit D

to Ji Declaration

Highly Confidential

Exhibit E

to Ji Declaration

Highly Confidential

Exhibit F

to Ji Declaration

Highly Confidential

